

Designing an Efficient Multimodal Biometric System using Palmprint and Speech Signal

Mahesh P.K.¹, M.N. Shanmukha Swamy²

¹ JSS Research Foundation, SJCE, Mysore, India

Email: mahesh24pk@gmail.com

² JSS Research Foundation, SJCE, Mysore, India

Email: mnsjce@gmail.com

Abstract— This paper proposes a multimodal biometric system using palmprint and speech signal. In this paper, we propose a novel approaches for both the modalities. We extract the features using Subband Cepstral Coefficients for speech signal and Modified Canonical method for palmprint. The individual feature score are passed to the fusion level. Also we have proposed a new fusion method called weighted score. This system is tested on clean and degraded database collected by the authors for more than 300 subjects. The results show significant improvement in the recognition rate.

Index Terms—Multimodal biometrics, Speech signal, Palmprint, Fusion

I. INTRODUCTION

A unimodal biometric authentication, which identifies an individual person using physiological and/or behavioral characteristics, such as palmprint, face, fingerprints, hand geometry, iris, retina, vein and speech. These methods are more reliable and capable than knowledge-based (e.g. Password) or token-based (e.g. Key) techniques. Since biometric features are hardly stolen or forgotten.

However, a single biometric feature sometimes fails to be exact enough for verifying the identity of a person. By combining multiple modalities enhanced performance reliability could be achieved. Due to its promising applications as well as the theoretical challenges, multimodal biometric has drawn more and more attention in recent years [1]. Speech Signal and palmprint multimodal biometrics are advantageous due to the use of non-invasive and low-cost speech and image acquisition. In this method we can easily acquire palmprint images using digital cameras, touchless sensors and speech signal using microphone. Existing studies in this approach [2, 3] employ holistic features for palmprint and speech signal representation and results are shown with different techniques of fusion and algorithms.

Multimodal system also provides anti-spoofing measures by making it difficult for an intruder to spoof multiple biometric traits simultaneously. However, an integration scheme is required to fuse the information presented by the individual modalities.

This paper presents a novel fusion strategy for personal identification using speech signal and palmprint features at the features level fusion Scheme. The proposed paper shows that integration of speech signal and palmprint biometrics can achieve higher performance that may not be possible using a single biometric indicator alone. We extract the

features using modified canonical form method for palmprint and Subband Cepstral Coefficients for speech. Integrating these two features at fusion level, which gives better performance and better accuracy. Which gives better performance and better accuracy for both traits (speech signal & palmprint).

The rest of this paper is organized as follows. Section 2 presents the system structure, which is used to increase the performance of individual biometric trait; multiple classifiers are combined using matching scores. Section 3 presents feature extraction method used for speech signal and section 4 for palmprint. Section 5, the individual traits are fused at matching score level based on weighted sum of score technique. Finally, the experimental results are given in section 6. Conclusions are given in the last section.

II. SYSTEM OVERVIEW

The block diagram of a multimodal biometric system using two (palm and speech) modalities for human recognition system is shown in Figure 1. It consists of three main blocks, that of Preprocessing, Feature extraction and Fusion. Preprocessing and feature extraction are performed in parallel for the two modalities. The preprocessing of the audio signal under noisy conditions includes signal enhancement, tracking environment and channel noise, feature estimation and smoothing [4]. The preprocessing of the palmprint typically consists of the challenging problems of detecting and tracking of the palm and the important palm features.

Further, features are extracted from the training and testing images and speech signal respectively, and then matched to find the similarity between two feature sets. The matching scores generated from the individual recognizers are passed to the decision module where a person is declared as genuine or an imposter.

III. SUBBAND BASED CEPSTRAL COEFFICIENTS AND GAUSSIAN MIXTURE MODEL

A. Subband Decomposition via Wavelet Packets

A detailed discussion of wavelet analysis is beyond the scope of this paper and we therefore refer interested readers to a more complete discussion presented in [5]. In continuous time, the Wavelet Transform is defined as the inner product of a signal $x(t)$ with a collection of wavelet functions $y_{ab}(t)$ in which the wavelet functions are scaled (by a) and translated

(by b) versions of the prototype wavelet $y(t)$.

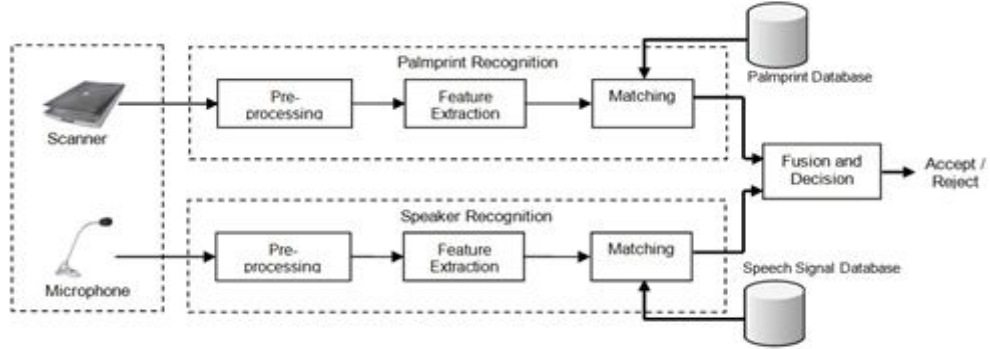


Figure 1. Block diagram of the proposed multimodal biometric verification system

$$\psi_{a,b}(t) = \psi\left(\frac{t-b}{a}\right) dt \quad (1)$$

$$W_{\psi} x(a,b) = \frac{1}{\sqrt{a}} \int_{-\infty}^{+\infty} x(t) \psi\left(\frac{t-b}{a}\right) dt \quad (2)$$

Discrete time implementation of wavelets and wavelet packets are based on the iteration of two channel filter banks which are subject to certain constraints, such as low pass and/or high pass branches on each level followed by a sub sampling-by-two unit. Unlike the wavelet transform which is obtained by iterating on the low pass branch, the filterbank tree can be iterated on either branch at any level, resulting in a tree structured filterbank which we call a wavelet packet filterbank tree. The resultant transform creates a division of the frequency domain that represents the signal optimally with respect to the applied metric while allowing perfect reconstruction of the original signal. Because of the nature of the analysis in the frequency domain it is also called subband decomposition where subbands are determined by a wavelet packet filterbank tree.

B. Wavelet Packet Transform Based Feature Extraction Procedure

Here, speech is assumed to be sampled at 8 kHz. A frame size of 24msec with a 10msec skip rate is used to derive the Subband based Cepstral Coefficients features, whereas a 20msec frame with the same skip rate is used to derive the MFCCs. We have used the same configuration proposed in [6] for MFCC. Next, the speech frame is Hamming windowed and pre-emphasized.

The proposed tree assigns more subbands between low to mid frequencies while keeping roughly a log-like distribution of the subbands across frequency. The wavelet packet transform is computed for the given wavelet tree, which results in a sequence of subband signals or equivalently the wavelet packet transform coefficients, at the leaves of the Tree. In effect, each of these subband signals contains only restricted frequency information due to inherent bandpass filtering. The wavelet packet tree is given in Figure 2. The energy of the sub-signals for each subband is computed and then scaled by the number of transform coefficients in that subband.

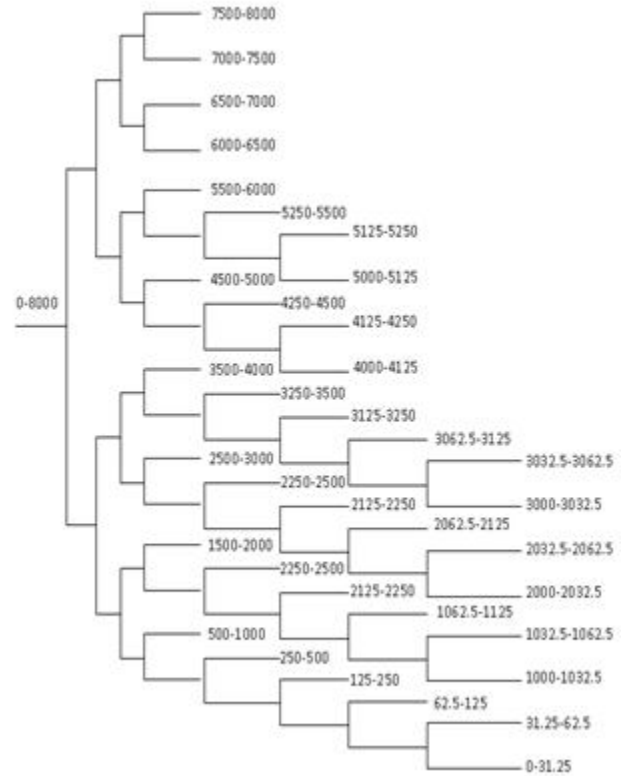


Figure 2. Wavelet Packet Tree

The subband signal energies are computed for each frame as,

$$S_i = \frac{\sum_{mel} [(W_{\psi})(i), m]}{N_i} \quad (3)$$

W_{ψ} : Wavelet packet transform of signal x,
 i : subband frequency index ($i=1,2,...L$),
 N_i : number of coefficients in the i^{th} subband.

C. Subband based Cepstral Coefficients

As in MFCCs the derivation of coefficients is performed in two stages. The first stage is the computation filterbank energies and the second stage would be the decorrelation of the log filterbank energies with a DCT to obtain the MFCC. The derivation of the Subband Based Cepstral coefficients follows the same process except that the filterbank energies

are derived using the wavelet packet transform rather than the short-time Fourier transform. It will be shown that these features outperform MFCCs. We attribute this to the computation of subband signals with smooth filters. The effect of filtering as a result of tracing through the low-pass/high-pass branches of the wavelet packet tree, is much smoother due to the balance in time-frequency representation. We believe that this will contribute to improved speech/speaker characterization over MFCC. Subband Based Cepstral coefficients are derived from subband energies by applying the Discrete Cosine Transformation:

$$SBC(n) = \sum_{i=1}^L \log S_i \cos\left(\frac{n(i-0.5)}{L} \pi\right), n = 1, \dots, n' \quad (4)$$

where n' is the number of SBC coefficients and L is the total number of frequency bands. Because of the similarity to root-cepstral [7] analysis, they are termed as subband based cepstral coefficients.

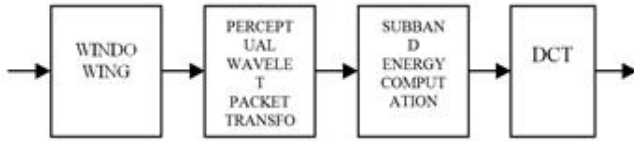


Figure 3. Block diagram for Wavelet Packet Transform based feature extraction procedure

D. The Gaussian Mixture Model

In this study, a Gaussian Mixture Model approach proposed in [8] is used where speakers are modeled as a mixture of Gaussian densities. The use of this model is motivated by the interpretation that the Gaussian components represent some general speaker-dependent spectral shapes and the capability of Gaussian mixtures to model arbitrary densities.

The Gaussian Mixture Model is a linear combination of M Gaussian mixture densities, and given by the equation,

$$p(\vec{x} | \lambda) = \sum_{i=1}^M p_i b_i(\vec{x}) \quad (5)$$

Where \vec{x} is a D -dimensional random vector, $b_i(\vec{x})$, $i=1, \dots, M$ are the component densities and p_i , $i=1, \dots, M$ are the mixture weights. Each component density is a D -dimensional Gaussian function of the form

$$b_i(\vec{x}) = \frac{1}{(2\pi)^{D/2} |\Sigma_i|^{1/2}} \exp\left\{-\frac{1}{2}(\vec{x} - \vec{\mu})^T \Sigma_i^{-1}(\vec{x} - \vec{\mu})\right\} \quad (6)$$

Where $\vec{\mu}$ denotes the mean vector and Σ_i denotes the covariance matrix. The mixture weights satisfy the law of total

probability, $\sum_{i=1}^M p_i = 1$. The major advantage of this representation of speaker models is the mathematical

tractability where the complete Gaussian mixture density is represented by only the mean vectors, covariance matrices and mixture weights from all component densities.

IV. FEATURE EXTRACTION USING MODIFIED CANONICAL FORM METHOD

Features are the attributes or values extracted to get the unique characteristics from the image and speech signal.

A. Palmprint feature extraction methodology

Details of the algorithm are as follows:

1) Identify hand image from background

Our designed system is such that palmprint images are captured using contact-less without pegs, keeping the image background relatively uniform and relatively low intensity when compared to the hand image. Using the statistical information of the background, the algorithm estimates an adaptive threshold to segment the image of the hand from the background. Pixels with intensity above the threshold are considered to be part of the hand image.

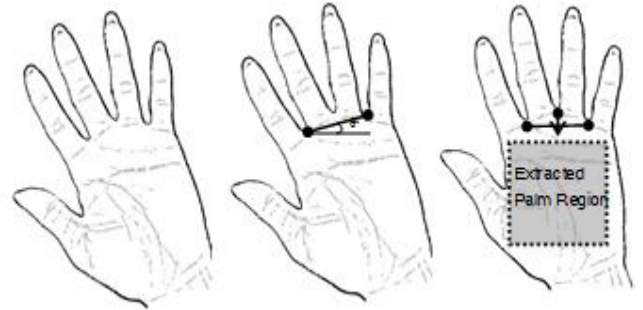


Figure 4. Schematic diagram of image alignment



Figure 5. Segmentation of ROI

2) Locate region-of-interest

The palm area is extracted from the binary image of the hand. After translating the original image into binary image, we find two key positioning points in the palmprint image using automatic detecting method. The first valley in the graph is the gaps between little finger and ring finger, Key Point 1. The third valley in the graph is the gaps between middle finger and index finger, Key Point 2. The key point is circled in Figure 4. The hand image is rotated by θ degrees. The hand images are rotated to align the hand images into a predefined direction. θ is calculated using the key points as shown in the Figure 4. Since the size of the original image is large, a smaller hand image is cropped out from the original hand image after image alignment using key points. Figure 5

shows the proposed image alignment and ROI selection method.

B. Modified Canonical Form Method

The “Eigenpalm” method proposed by Turk and Pentland [9][10] is based on Karhunen-Loeve Expression and we are motivated by this work for efficiently representing picture of images. The Eigen method presented by Turk and Pentland finds the principal components (Karhunen-Loeve Expression) of the image distribution or the eigenvectors of the covariance matrix of the set of images. These eigenvectors can be thought as set of features, which together characterized between images

Let an image $I(x, y)$ be a two dimensional array of intensity values or a vector of dimension n . Let the training set of images be $I_1, I_2, I_3, \dots, I_n$. The average image of the set is defined by

$$\Psi = \frac{1}{N} \sum_{i=1}^n I_i \quad (7)$$

Each image differed from the average by the vector.

$$\phi_i = I_i - \Psi \quad (8)$$

This set of very large vectors is subjected to principal component analysis which seeks a set of K orthonormal vectors $V_k, K=1, \dots, K$ and their associated eigenvalues λ_k which best describe the distribution of data. The vectors V_k and scalars λ_k are the eigenvectors and eigenvalues of the covariance matrix:

$$C = \frac{1}{N} \sum_{i=1}^N \phi_i \phi_i^T = A A^T \quad (9)$$

Where the matrix $A = [\phi_1, \phi_2, \dots, \phi_N]$ finding the eigenvectors of matrix $C_{n \times n}$ is computationally intensive. However, the eigenvectors of C can determine by first finding the eigenvectors of much smaller matrix of size $N \times N$ and taking a linear combination of the resulting vectors [4].

The modified canonical method proposed in this paper is based on Eigen values and Eigen vectors. These Eigen values can be thought a set of features which together characterized between images.

Let \hat{P} be the normalized modal matrix of I , the diagonal matrix is given by

$$D = \hat{P}^{-1} C \hat{P} \quad (10)$$

Where $\hat{P} = \frac{V_{k_{ij}}}{X_i}$ and

$$X_i = \sqrt{\sum V_{k_{ij}}^2}, i, j=1, 2, 3, \dots, n \quad (11)$$

Then the quadratic form Q is given by

$$Q = X^T D X \quad (12)$$

The following steps are considered for the feature extraction:

- Select the palm image for the input
- Pre-process the image
- Determine the eigen values and eigen vectors of the image
- Use the canonical form for the feature extraction.

C. Euclidean Distance

Let an arbitrary instance X be described by the feature vector

$$X = [a_1(x), a_2(x), \dots, a_n(x)] \quad (13)$$

Where $a_r(x)$ denotes the value of the r^{th} attribute of instance x . Then the distance between two instances x_i and x_j is defined to be $d(x_i, x_j)$;

$$d(x_i, x_j) = \sqrt{\sum_{r=1}^n (a_r(X_i) - a_r(X_j))^2} \quad (14)$$

D. Score Normalization

This step brings both matching scores between 0 and 1 [11]. The normalization of both the scores are done by

$$N_{\text{Speech}} = \frac{MS_{\text{Speech}} - \min_{\text{Speech}}}{\max_{\text{Speech}} - \min_{\text{Speech}}} \quad (15)$$

$$N_{\text{Palm}} = \frac{MS_{\text{Palm}} - \min_{\text{Palm}}}{\max_{\text{Palm}} - \min_{\text{Palm}}} \quad (16)$$

Where \min_{Speech} and \max_{Speech} are the minimum and maximum scores for speech signal recognition and $\min_{\text{Palmprint}}$ and $\max_{\text{Palmprint}}$ are the corresponding values obtained from palmprint trait.

E. Generation of Similarity Scores

Note that the normalized score of palmprint which is obtained through Haar Wavelet gives the information of dissimilarity between the feature vectors of two given images while the normalized score from speech signal gives a similarity measure. So to fuse both the score, there is a need to make both the scores as either similarity or dissimilarity measure. In this paper, the normalized score of palmprint is converted to similarity measure by

$$N'_{\text{Palm}} = 1 - N_{\text{Palm}} \quad (17)$$

V. FUSION

The biometrics systems are integrated at multi-modality level to improve the performance of the verification system. At multi-modality level, matching score are combined to give a final score. The following steps are performed for fusion:

1. Given a query image and speech signal as input, features are extracted by the individual recognition and then the matching score of each individual trait is calculated.
2. The weights a and b are calculated using FAR and FRR.
3. Finally, the final score after combining the matching score

of each trait is calculated by weighted sum of score technique.

$$MS_{fusion} = \frac{a * MS_{Palm} + b * MS_{Speech}}{2} \quad (18)$$

Where a and b are the weights assigned to both the traits. The final matching score (MS_{fusion}) is compared against a certain threshold value to recognize the person as genuine or an imposter.

VI. EXPERIMENTAL RESULTS

This section shows the experimental results of our approach with Modified Canonical method and Subband based Cepstral coefficients for palmprint and Speech respectively. We evaluate the proposed multimodal system on a data set including more than 300 subjects taking 6 different samples, also we have experimented with two different conditions (Cleaned and Degraded data). The training database contains a palmprint images and speech signal for each individual for each subject.

The comparison of both unimodal systems (palm and speech modality) and a bimodal system is given in Table 1 & 2. It can be seen that the fusion of palmprint and speech features improves the verification score. The experiments show that EER is reduced to 3.54% in clean database and 9.17% in degraded database.

TABLE I. THE FAR AND FRR OF PALMPRINT AND SPEECH SIGNAL IN CLEAN AND DEGRADED CONDITIONS

| Modality | Method | Database | Classifier | FAR% | FRR% |
|---------------|--------|----------------------------------|------------|-------|-------|
| Speech Signal | SBC | Clean database of 300 subjects | GMM | 03.54 | 10.73 |
| Palmprint | MCF | | --- | 06.73 | 15.24 |
| Speech Signal | SBC | Degraded database of 50 subjects | GMM | 46.67 | 31.67 |
| Palmprint | MCF | | --- | 23.33 | 26.67 |

TABLE II. THE FAR AND FRR AFTER FUSION

| Methods | | Database | FAR% | FRR% |
|-----------|---------------|----------------------------------|-------|-------|
| Palmprint | Speech Signal | | | |
| MCF | SBC | Clean database of 300 subjects | 02.76 | 04.32 |
| MCF | SBC | Degraded database of 50 subjects | 06.67 | 11.67 |

CONCLUSIONS

Biometric systems are widely used to overcome the traditional methods of authentication. But the unimodal biometric system fails in case of biometric data for particular trait. This paper proposes a new method in selecting and dividing the ROI for analysis of palmprint. The new method utilizes the maximum palm region of a person to attain feature extraction. More importantly, it can cope with slight variations, in terms of rotation, translation, and size difference, in images captured from the same person. The experimental results show that the performance of palmprint-based unimodal system and speech-based unimodal system fails to meet the requirement. Fusion at the matching-score level is used to improve the performance of the system. The psychological effects of such multimodal system should also not be disregarded and it is likely that a system using multiple modalities would seem harder to cheat to any potential impostors.

In the future we plan to test whether setting the user specific weights to different modalities can be used to improve a system's performance.

REFERENCES

- [1] A. A. Ross, K. Nandakumar, and A. K. Jain. Handbook of Multibiometrics. Springer-Verlag, 2006.
- [2] Mahesh P.K. and M.N. Shanmukhaswamy. Comprehensive Framework to Human Recognition Using Palmprint and Speech Signal. In *Springer-Verlag Berlin Heidelberg 2011*.
- [3] Mahesh P.K. and M.N. Shanmukhaswamy. Integration of multiple cues for human authentication system. In *Procedia Computer Science*, Volume 2, 2010, Pages 188-194.
- [4] Jr., J. D., Hansen, J., and Proakis, J. Discrete Time Processing of Speech Signals, *second ed. IEEE Press, New York, 2000*.
- [5] O. Rioul and M. Vetterli, "Wavelets and Signal Processing," *IEEE Signal Proc. Magazine*, vol. 8(4), pp. 11-38, 1991.
- [6] D. A. Reynolds and R. C. Rose, "Robust Text Independent Speaker Identification Using Gaussian Mixture Speaker Models" *IEEE Transactions on SAP*, vol.3, pp. 72-83, 1995.
- [7] P. Alexandre and P. Lockwood, "Root cepstral analysis: A unified view: Application to speech processing in car noise environments," *Speech Communication*, v.12, pp. 277-288,1993.
- [8] D. A. Reynolds, "Experimental Evaluation of Features for Robust Speaker Identification," *IEEE Transactions on SAP*, vol. 2. Pp. 639-643,1994.
- [9] Turk and A. Pentland, "Face Recognition using Eigenfaces", in *Proceeding of International Conference on Pattern Recognition*, pp. 591-1991.
- [10] Turk and A. Pentland, "Face Recognition using Eigenfaces", *Journals of Cognitive Neuroscience*, March 1991.
- [11] A. K. Jain, K. Nandakumar, & A. Ross, Score Normalization in multimodal biometric systems. *The Journal of Pattern Recognition Society*, 38(12), 2005, 2270-2285.